

AD-A080 601

WISCONSIN UNIV-MADISON DEPT OF STATISTICS

F/G 12/1

TRANSFORMATION OF GROUPED OR CENSORED DATA TO NEAR NORMALITY.(U)

NOV 79 V M GUERRERO, R A JOHNSON

N00014-78-C-0722

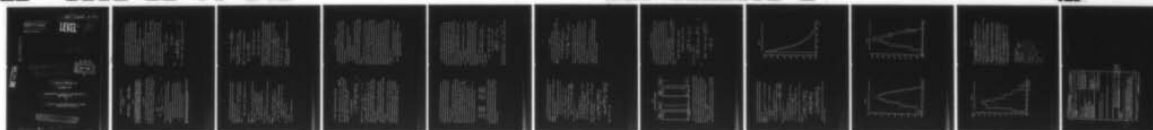
UNCLASSIFIED

UNIS-DS-79-542

AD-15660.2-M

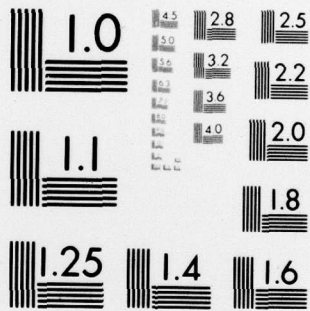
MI

1 OF 1  
AD  
A080601



END  
DATE  
FILMED

3 - 80  
DDC



MICROCOPY RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS-1963-A

DEPARTMENT OF STATISTICS  
UNIVERSITY OF WISCONSIN

Madison, Wisconsin

18 ARO 15660.2-M  
**LEVEL** # ①

ADA080601

DDC FILE COPY

14 UWIS-DS-79-542

15 N00014-78-C-0722

DDC  
RECEIVED  
FEB 13 1980  
A

9 TECHNICAL REPORT NO. 542

11 November 1979

6 TRANSFORMATION OF GROUPED OR CENSORED DATA TO  
NEAR NORMALITY

by

10 Victor M. Guerrero and Richard A. Johnson  
University of Wisconsin

THE VIEW, OPINIONS, AND/OR FINDINGS CONTAINED IN THIS REPORT  
ARE THOSE OF THE AUTHOR(S) AND SHOULD NOT BE CONSTRUED AS  
AN OFFICIAL DEPARTMENT OF THE ARMY POSITION, POLICY, OR DE-  
CISION, UNLESS SO DESIGNATED BY OTHER DOCUMENTATION.

✓  
400 243

80

2

8

106

JOB

# Transformation of Grouped or Censored Data to

## Near Normality

by

Victor M. Guerrero  
University of Wisconsin

Richard A. Johnson  
University of Wisconsin

## Summary

Box and Cox (1964) proposed a power transformation which has proven utility for transforming ungrouped data to near normality. In this paper, we extend its applicability to two frequently encountered situations: (i) grouped data and (ii) Type I censored data. The asymptotic properties of the estimators are derived and illustrative examples are presented.

## 1. Introduction

Let us consider a random sample  $X_1, \dots, X_n$  from an absolutely continuous distribution with probability density function (pdf)  $g$  concentrated on  $(0, \infty)$ . We assume the existence of a unique vector of parameter values  $\theta_0 = (\nu_0, \sigma_0^2, \lambda_0)$  such that the distribution of the transformed variables  $X_j^{(\lambda)}$  defined by

$$X_j^{(\lambda)} = \begin{cases} \frac{X_j^{\lambda_0-1}}{\lambda_0} & \text{if } \lambda_0 \neq 0 \\ \log(X_j) & \text{if } \lambda_0 = 0 \end{cases} \quad \text{for } j=1, \dots, n \quad (1)$$

is "closest" to a normal distribution with parameters  $\nu_0$  and  $\sigma_0^2$ .

The problem to be considered here is that of obtaining the Maximum Likelihood Estimator (MLE) of  $\theta$  when:

- (i) the original random variables are unobserved and the only available information is the number of observations falling within arbitrary, but specified, intervals of the real line, or
- (ii) some exact observations are available besides the counts of observations in other specified intervals.

Usually, when dealing with incomplete data, explicit expressions for the MLE's are not available. However, it is possible to gain some insight into the proposed procedure by studying their asymptotic properties as in Hernández and Johnson (1979). Thus we establish strong consistency and asymptotic normality of the MLE and identify its limit with a minimum Kullback-Leibler information number property.

Examples of both grouped and censored data are considered.

## 2. Grouped Data

Let the sample be grouped into  $k(k \geq 3)$  prespecified intervals denoted by  $D_1 = [a_0, a_1), D_2 = [a_1, a_2), \dots, D_k = [a_{k-1}, a_k)$  where  $0 = a_0 < a_1 < \dots < a_{k-1} < a_k = \infty$ . The count of the number of observations in interval  $D_i$  will be denoted by  $n_i$ , and the total sample size is  $n = \sum_{i=1}^k n_i$ . In order to obtain the MLE of  $\theta$ , we tentatively assume (pretend),

$X_j^{(\lambda)} \sim N(\nu_0, \sigma_0^2)$  for  $j = 1, \dots, n$ . Because  $X_j$  is positive,  $X_j^{(\lambda)}$  cannot have a normal distribution, except possibly for  $\lambda_0 = 0$ . Under this assumption the log-likelihood becomes

$$L_n(\theta) = \log(n!) - \sum_{i=1}^k \log(n_i!) + \sum_{i=1}^k n_i \log[p_i(\theta)] \quad (2)$$

where

$$p_i(\theta) = \phi\left(\frac{\lambda_0}{\sigma_0} \frac{a_i - \mu}{\sigma_0}\right) - \phi\left(\frac{\lambda_0}{\sigma_0} \frac{a_{i-1} - \mu}{\sigma_0}\right) \quad \text{for } i = 1, \dots, k \quad (3)$$

$\phi(\lambda) = -\infty$  for every  $\lambda$  and  $\phi(x) = \int_{-\infty}^x \phi(t) dt$  with  $\phi(t) = (2\pi)^{-1/2} e^{-t^2/2}$ .

\*This research was sponsored by the Office of Naval Research under Grant No. N00014-78-C-0722. (Also funded by Army Research Office).

We state without proofs two lemmas which are used to establish (1) strong consistency and (2) asymptotic normality of the MLE  $\hat{\theta}_n$  (see Guerrero (1979) for the proofs.)

Lemma 1. Let  $D_1 = [a_{1,1}, a_{1,1}]$ ,  $q_1 = \int_{D_1} g(x)dx$ ,  $\hat{p}_{1,n} = n_1/n$  and  $p_1(\theta)$  be as in (3) for  $i = 1, \dots, k$ . For  $\Omega$  a compact subset of  $\mathbb{R}^3$  we have

$$\left| \sum_{i=1}^k \hat{p}_{1,n} \log[p_1(\theta)] - \sum_{i=1}^k q_1 \log[p_1(\theta)] \right| \xrightarrow{a.s.} 0 \text{ as } n \rightarrow \infty$$

uniformly in  $\theta \in \Omega$ .

Lemma 2. For fixed  $n$ , let  $s_n(\theta)$  be the log-likelihood function (2) and  $\hat{\theta}_n$  the MLE of  $\theta$ . For  $\Omega$  a compact set, assume

- (i)  $f$  is a continuous function from  $\Omega$  to  $\mathbb{R}$ ,
- (ii)  $\frac{1}{n} s_n(\theta) \xrightarrow{a.s.} f(\theta)$  uniformly in  $\theta \in \Omega$ ,
- (iii)  $f(\theta)$  has a unique global maximum at  $\theta = \theta_0$ .

Then,  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$  as  $n \rightarrow \infty$ .

We now proceed to state and sketch the proof of the main result of this section.

Theorem 3. Let  $q_i$  and  $p_i(\theta)$  be as in Lemma 1. If

- (i) the parameter space is a compact subset of  $\mathbb{R}^3$

$$(ii) H(\theta) = \sum_{i=1}^k q_i \log \left[ \frac{p_i(\theta)}{q_i} \right] \text{ has a unique global maximum}$$

as a function of  $\theta = (\theta_1, \theta_2, \theta_3)' = (\mu, \sigma, \lambda)'$ , and this

is attained at  $\theta = \theta_0$ .

Then, (i)  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$  as  $n \rightarrow \infty$ .

If further,

- (iii)  $\theta_0$  is an interior point of  $\Omega$ ,

(iv) the Hessian of  $H(\theta)$ ,  $\nabla^2 H(\theta) = \left( \frac{\partial^2 H(\theta)}{\partial \theta_u \partial \theta_v} \right)_{3 \times 3}$  is nonsingular at  $\theta_0$ .

Then, (2)  $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_3(0, VV')$  as  $n \rightarrow \infty$ , where  $V = [\nabla^2 H(\theta_0)]^{-1}$  and  $V = (v_{uv})_{3 \times 3}$  is given by (4) below.

Proof: Let  $\hat{p}_{1,n} = \frac{n_1}{n}$ . Stirling's approximation yields

$$\begin{aligned} \left| \frac{1}{n} s_n(\theta) - \sum_{i=1}^k q_i \log \left[ \frac{p_i(\theta)}{q_i} \right] \right| &\leq \left| \sum_{i=1}^k \hat{p}_{1,n} \log[p_1(\theta)] - \sum_{i=1}^k q_i \log[p_1(\theta)] \right| \\ &+ \left| \sum_{i=1}^k \hat{p}_{1,n} \log(\hat{p}_{1,n}) - \sum_{i=1}^k q_i \log(q_i) \right| + o(1) \end{aligned}$$

Because of Lemma 1, the right hand side converges to zero, with probability one, uniformly in  $\theta \in \Omega$ . The consistency of  $\hat{\theta}_n$  follows from Lemma 2.

To establish asymptotic normality, we obtain the gradient and the Hessian of the log-likelihood function by straightforward differentiation.

The Hessian of  $s_n(\theta)$ , with elements  $\frac{\partial^2 s_n(\theta)}{\partial \theta_u \partial \theta_v}$  for  $u, v=1, 2, 3$  is readily seen to be continuous on  $\Omega$ . So, using Taylor's formula,

$$\frac{1}{\sqrt{n}} \nabla s_n(\hat{\theta}) = \frac{1}{\sqrt{n}} \nabla s_n(\theta_0) + \frac{1}{n} \nabla^2 s_n(\theta_0) [\sqrt{n}(\hat{\theta}_n - \theta_0)]$$

with  $\theta_0 = \gamma_n \theta_0 + (1 - \gamma_n) \hat{\theta}_n$ ,  $0 < \gamma_n < 1$ .

We conclude that  $\frac{1}{\sqrt{n}} \nabla s_n(\theta_0)$  and  $-\frac{1}{n} \nabla^2 s_n(\theta_0) [\sqrt{n}(\hat{\theta}_n - \theta_0)]$  have the same limiting distribution. Write



$$\frac{1}{n} \sum_{j=1}^n v_{ij}(\theta_0) = \frac{1}{n} \sum_{j=1}^n \left[ \sum_{i=1}^k I_{D_i}(x_j) \alpha_{ij}(\theta_0) \right] = \frac{1}{n} \sum_{j=1}^n z_j(\theta_0)$$

where  $\alpha_{ij}(\theta_0) = (\alpha_{11}(\theta_0), \alpha_{12}(\theta_0), \alpha_{13}(\theta_0))'$  with  $\alpha_{1r}(\theta_0) = \frac{\partial \log p_1(\theta)}{\partial \theta_r} \bigg|_{\theta_0}$  for  $r = 1, 2, 3$ . Asymptotic normality follows from the multivariate Central Limit Theorem since the random vectors  $z_1(\theta_0), \dots, z_n(\theta_0)$  are iid with  $E[z_1(\theta_0)] = \nabla H(\theta_0) = \bar{0}$  and have covariance elements

$$\begin{aligned} v_{uv} &= \sum_{i=1}^k q_i \alpha_{iu}(\theta_0) \alpha_{iv}(\theta_0) \\ &= \sum_{i=1}^k q_i \left( \frac{\partial \log p_1(\theta)}{\partial \theta_u} \bigg|_{\theta_0} \right) \left( \frac{\partial \log p_1(\theta)}{\partial \theta_v} \bigg|_{\theta_0} \right). \end{aligned} \quad (4)$$

Remark: Setting  $Q = (q_i | q_i = \int_{D_i} g(x) dx > 0, i = 1, \dots, k)$  and

$$M = (p_i(\theta) | p_i(\theta) = \phi \left( \frac{\frac{\partial(\lambda)}{\partial \lambda} - \mu}{\sigma} \right) - \phi \left( \frac{\frac{\partial(\lambda)}{\partial \lambda} - \mu}{\sigma} \right), i = 1, \dots, k), \text{ it}$$

follows easily that, under the conditions of Theorem 3,  $\theta_0$  is also that value of  $\theta$  which minimizes the Kullback-Leibler information number  $I[Q, M]$ . A similar interpretation for ungrouped data is given by Hernández and Johnson (1979).

Hence, obtaining the MLE of  $\theta$  using the (wrong) log-likelihood function (2) is asymptotically equivalent to finding the minimum of the Kullback-Leibler information number between the true probability distribution  $Q$  and a normal.

The importance of the true pdf  $g$  is reflected in the asymptotic variance of  $\hat{\theta}_n$  derived in Theorem 3. If one is faced with the task of transforming a grouped sample to near normality, the true probabilities  $q_i$  can be estimated by the observed frequencies  $\hat{p}_{i,n}$ . Doing this, one obtains a consistent estimate of the asymptotic variance of  $\hat{\theta}_n$ .

### 3. Application

Because grouping results in loss of information, it is preferable to transform the original (ungrouped) data. Only when the original data are not available should grouped observations be transformed. We may want to transform to normality

- (a) for better description or for model building purposes,
- (b) to interpolate between the given data,
- (c) to smooth count data.

The technique we propose will probably be most appropriate when the sample size is moderate so that the large sample theory gives some assurances, yet not so large that one would use splines or other nonparametric smoothing procedures.

In order to find the MLE of  $\theta = (\mu, \sigma, \lambda)'$  in our transformed normal approach, we need to use an iterative procedure and this requires some initial estimates. Several methods have been devised for obtaining approximate MLE's for  $\mu$  and  $\sigma$  (c.f. Lindley (1950), Bunn and Sidebottom (1976)) which may be used for getting initial values for those parameters. But the problem still remains with respect to selecting an initial value for  $\lambda$ .

In practice, we can apply a specialization to our case, of a two-stage procedure proposed by Richards (1961) (the assumptions under which this procedure is valid were wrongly stated by Richards, but later Kale (1963) validated the method by assuming that the conditions of the Implicit Function Theorem hold). This method enables us to obtain the value  $\hat{\theta}_n = (\hat{\mu}_n, \hat{\sigma}_n, \hat{\lambda}_n)'$ , for  $n$  fixed, as follows:

- Step (1) fix  $\lambda$  and maximize  $L_n(\theta)$  with respect to  $\mu$  and  $\sigma$ , thus obtaining a value  $\hat{\theta}_n(\lambda) = (\hat{\mu}_n(\lambda), \hat{\sigma}_n(\lambda), \lambda)'$ .
- Step (2) maximize the already partially maximized log-likelihood  $L_n(\hat{\theta}_n(\lambda))$  with respect to  $\lambda$  to get  $\hat{\theta}_n$ .

In practice it is not possible to get explicit expressions for

$\hat{\mu}_n(\lambda)$  and  $\hat{\sigma}_n(\lambda)$ , therefore the procedure relies heavily upon the use of the Implicit Function Theorem to assure the existence of such quantities. To find the values  $\hat{\mu}_n(\lambda)$  we solve, iteratively, the maximum likelihood equations for  $\mu$  and  $\sigma$ .

We illustrate our procedure for the life length distribution in Figure 1. The data are the grouped ages of the Mexican population in 1966 which appear in Keyfitz and Flieger (1971 p. 344). Because of the large sample size, other smoothing techniques may be more appropriate. The value of  $\lambda$  which maximized the log-likelihood function is  $\hat{\lambda} = .34995$  and for this value we obtained  $\hat{\mu} = 4.775151$  and  $\hat{\sigma} = 2.712447$ . These values were used to draw the histogram of the transformed data with an overlay of the appropriate normal pdf in Figure 2. Shown in Figure 1 is the corresponding transformed normal pdf. From the estimated variance-covariance matrix of  $\hat{\theta}$  we have

$$\begin{aligned}\text{Var}(\hat{\mu}) &= .0000000726 & \text{Var}(\hat{\sigma}) &= .0000000634 \\ \text{Var}(\hat{\lambda}) &= .0000000011 & \text{Cov}(\hat{\mu}, \hat{\sigma}) &= .0000000576 \\ \text{Cov}(\hat{\mu}, \hat{\lambda}) &= .0000000079 & \text{Cov}(\hat{\sigma}, \hat{\lambda}) &= .0000000078\end{aligned}$$

From Figures 1 and 2 we observe that strict normality was not achieved, although the transformation did yield a nearly symmetrical distribution. Draper and Cox (1969) noticed the fact that in some cases of ungrouped data, even when the transformation procedure does not yield normality, it helps to "regularize" data.

#### 4. Grouped Data Combined with Exact Observations in the Tails.

Here we consider a censoring scheme where exact values are recorded only for observations either in the lower or the upper tail, or for both extremes. The rest of the sample is grouped. More precisely, the values of all observations below the (specified) fixed value  $a_1$  or above the fixed value  $a_{k-1}$  are recorded exactly, while those in between are grouped into the intervals  $[a_1, a_2), \dots, [a_{k-2}, a_{k-1})$ , where  $0 < a_1 < \dots < a_{k-1} < \infty$ . This is a form of fixed time censoring which produces a mixed sample of exact and grouped data.

As before, we denote the intervals  $[0, a_1), [a_1, a_2), \dots, [a_{k-1}, \infty)$  by  $D_1, D_2, \dots, D_k$  respectively.

The motivation behind this formulation is that it is good statistical practice to report extreme values exactly for further scrutiny. This may be done either to check for outliers or to delineate unique features of the phenomenon under study, particularly those related to the tail behavior of the distributions.

Pretending that  $X_j \sim N(\mu_0, \sigma_0^2)$ ,  $j = 1, \dots, n$  for some

$\theta_0 = (\mu_0, \sigma_0^2, \lambda_0)'$ , the log-likelihood of the (mixed) sample, of size  $n = \sum_{i=1}^k n_i$ , becomes

$$\begin{aligned}l_n(\theta|x) &= \log(n!) - \sum_{i=1}^k \log(n_i!) + \sum_{i=2}^{k-1} n_i \log[p_i(\theta)] \\ &\quad - \frac{1}{2} (n_1 + n_k) \log(2\pi) - (n_1 + n_k) \log(\sigma) \\ &\quad - \frac{1}{2} \sum_{j=1}^{n_1, n_k} \left( \frac{x_j(\lambda) - \mu}{\sigma} \right)^2 + (\lambda-1) \sum_{j=1}^{n_1, n_k} \log(x_j),\end{aligned}\quad (5)$$

where  $p_i(\theta)$  is defined in (3) for  $i = 2, \dots, k-1$  and where  $\sum_{j=1}^{n_1, n_k}$  stands for the sum over the smallest  $n_1$  and the largest  $n_k$  order statistics.

Now, let  $z_1^{(u)}(\theta|x)$  denote the log-likelihood corresponding to one ungrouped observation. Then we obtain the following theorem.

**Theorem 4.** Let  $q_i$  and  $p_i(\theta)$  be defined as in Lemma 1 and suppose the following conditions are satisfied:

- (i) the parameter space  $\Omega$  is the compact subset of  $R^3$  defined by  $\Omega = \{\theta = (\mu, \sigma, \lambda) : |\mu| \leq M, s_1 \leq \sigma \leq s_2, a \leq \lambda \leq b$  for some  $0 < M, s_1, s_2, b < \infty$  and  $-\infty < a < 0\}$ .
- (ii) the moments  $E_g(x^{2a})$  and  $E_g(x^{2b})$  are finite.

- (iii)  $H(\theta) + E_g[I_{0,1} \mu_0(x) z_1^{(u)}(\theta|x)]$  has a unique global maximum at  $\theta = \theta_0$ , where  $H(\theta) = \sum_{i=2}^{k-1} q_i \log \left[ \frac{p_i(\theta)}{q_i} \right]$ .
- Then, (1)  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$  as  $n \rightarrow \infty$ .

Furthermore, if:

- (iv)  $\theta_0$  is an interior point of  $\Omega$ ,
- (v) both  $E_g[x^a \log(x)]^2$  and  $E_g[x^b \log(x)]^2$  are finite,
- (vi)  $W(\theta_0) + E_g[I_{0,1} \mu_0(x) \cdot v_1^{(u)}(\theta_0|x)] = 0$ ,
- (vii)  $V = (\sigma^2 H(\theta_0) + E_g[I_{0,1} \mu_0(x) \cdot v^2 z_1^{(u)}(\theta_0|x)])^{-1}$  exists.

Then, (2)  $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_3(0, WV)$  as  $n \rightarrow \infty$ , with the elements of  $W = (w_{uv})$ ,  $u, v = 1, 2, 3$ , given by

$$w_{uv} = \sum_{i=2}^{k-1} q_i \left( \frac{\partial \log p_i(\theta)}{\partial \theta_u} \Big|_{\theta_0} \right) \left( \frac{\partial \log p_i(\theta)}{\partial \theta_v} \Big|_{\theta_0} \right) + E_g \left[ I_{0,1} \mu_0(x) \left( \frac{\partial z_1^{(u)}(\theta|x)}{\partial \theta_u} \Big|_{\theta_0} \right) \left( \frac{\partial z_1^{(v)}(\theta|x)}{\partial \theta_v} \Big|_{\theta_0} \right) \right]$$

Proof: See Guerrero (1979).

### 5. An Example

Consider the situation in which the observations

below some fixed point  $a_j$  and above the fixed point  $a_{k-1}$  are known exactly and the rest of the sample is given in grouped form. Writing  $A_j$  for  $a_j^{(1)} - \mu$  the likelihood equations for  $\mu$  and  $\sigma$  are

$$\frac{\partial L}{\partial \mu} = \sum_{i=2}^{k-1} \frac{n_i}{\sigma} \left[ \frac{\phi(A_{i-1}) - \phi(A_i)}{p_i(\theta)} \right] + \frac{1}{\sigma^2} \left( \sum_{j=1}^h x_j^{(\lambda)} - h\mu \right)$$

$$\frac{\partial L}{\partial \sigma} = \sum_{i=2}^{k-1} \frac{n_i}{\sigma} \left[ \frac{A_{i-1} \phi(A_{i-1}) - A_i \phi(A_i)}{p_i(\theta)} \right] - \frac{h}{\sigma} + \frac{1}{\sigma^3} \sum_{j=1}^h (x_j^{(\lambda)} - \mu)^2$$

with  $n = n_1 + n_k$  the total number of exact observations. These equations are to be solved iteratively for fixed  $\lambda$  values.

Data approximating our assumptions come from the Hage Surveys, in which a category like "under 100" has a footnote giving counts for such a category in much smaller intervals. For the purpose of illustration we use the data on Weekly Earnings of Secretaries (in Hartford, CT) given in Table 1 below.

We will assume that observations outside the range 130-250 are given exactly so the five observations falling in [120, 130) are placed at 121, 123, 125, 127, and 129, and the last three observations will be taken as 255, 285 and 325.



Table 1  
WEEKLY EARNINGS OF SECRETARIES

Weekly Earnings of	No. of Secretaries	Weekly Earnings of	No. of Secretaries
120 and under 130	5	230 and under 240	24
130	45	240	4
140	125	250	1
150	126	260	-
160	141	270	-
170	100	280	1
180	69	290	-
190	70	300	-
200	48	310	-
210	43	320	-
220	11	330	1
230		Total	814

Source: Industry Wage Surveys, December 1976. Bulletin 1988 (Table 75).

A plot of the transformed histogram with the "closest" normal approximation is shown as Figure 3 below. It can be observed that even though exact normality has not been achieved, the transformation worked reasonably well. The improvement in normality can also be appreciated by fitting the inverted normal  $p_c$  to the original data as shown in Figure 4. The MLE's in this case were found to be  $\hat{\lambda} = -.99$  (approximately the reciprocal transformation),  $\hat{\mu} = 1.003809$  and  $\hat{\sigma} = .000891$ .

The chi-square statistic for a goodness-of-fit test is  $\chi^2 = 63.59$  which rejects normality because of the drops in intervals [180, 190) and [220, 230).

### 6. Type I Censoring Procedures

Consider data that occurs when a censoring scheme operates in such a way that observations falling between the specified fixed limits  $a_1$  and  $a_2$  are recorded exactly, while only the counts of values below  $a_1$  and above  $a_2$  ( $0 < a_1 < a_2 < \infty$ ) are available. This is the usual double Type I censoring scheme which, for our case, divides the positive real line into three intervals, namely  $D_1 = (n, a_1)$ ,  $D_2 = [a_1, a_2]$  and  $D_3 = [a_2, \infty)$ , and yields a sample of grouped and ungrouped data.

We again pretend that  $X_j^{(\lambda)} \sim N(\mu_0, \sigma_0^2)$ ,  $j = 1, \dots, n$  for some parameter value  $\theta_0 = (\mu_0, \sigma_0, \lambda_0)$ . The log-likelihood of the mixed sample of size  $n$  is then

$$\begin{aligned} \ell_n(\theta|x) = & \log(n!) - \sum_{i=1}^3 \log(n_i!) + n_1 \log[p_1(\theta)] \\ & + n_3 \log[p_3(\theta)] - \frac{1}{2} n_2 \log(2\pi) - n_2 \log(\sigma) \\ & - \frac{1}{2} \sum_{j=n_1+1}^{n-n_3} \left( \frac{x_j(\lambda) - \mu}{\sigma} \right)^2 + (\lambda-1) \sum_{j=n_1+1}^{n-n_3} \log(x_j). \end{aligned}$$

with  $n_1 + n_2 + n_3 = n$  and where

$$p_1(\theta) = \phi\left(\frac{a_1(\lambda) - \mu}{\sigma}\right) \quad \text{and} \quad p_3(\theta) = 1 - \phi\left(\frac{a_2(\lambda) - \mu}{\sigma}\right). \quad (6)$$

Let

$$\begin{aligned} \ell_{n_2}^{(u)}(\theta|x) = & -\frac{1}{2} n_2 \log(2\pi) - n_2 \log(\sigma) - \frac{1}{2} \sum_{j=n_1+1}^{n-n_3} \left( \frac{x_j(\lambda) - \mu}{\sigma} \right)^2 \\ & + (\lambda-1) \sum_{j=n_1+1}^{n-n_3} \log(x_j) \end{aligned}$$

be the part of the log-likelihood linked to the ungrouped data. Again consistency to the minimum information number solution and asymptotic normality follow.

**Theorem 5.** Let  $q_1 = \int_0^x g(x)dx$  for  $1 = 1, 2, 3$  and let  $p_1(\theta)$  and  $p_3(\theta)$  be as in (6). If

- (i) the parameter space  $\Omega$  is the compact set defined in

**Theorem 4,**

- (ii) the moments  $E_g(x^{2a})$  and  $E_g(x^{2b})$  are finite,  
 (iii)  $H(\theta) + E_g[I_{D_2}(x)z_1^{(u)}(\theta|x)]$  has a unique global maximum

$$\text{at } \theta = \theta_0, \text{ where } H(\theta) = q_1 \log \left[ \frac{p_1(\theta)}{q_1} \right] + q_3 \log \left[ \frac{p_3(\theta)}{q_3} \right].$$

Then, (1)  $\hat{\theta}_n \xrightarrow{d.s.} \theta_0$  as  $n \rightarrow \infty$ .

Furthermore, if:

- (iv)  $\theta_0$  is an interior point of  $\Omega$ ,  
 (v) both  $E_g[x^b \log(x)]^2$  and  $E_g[x^b \log(x)]^2$  are finite,  
 (vi)  $\nabla H(\theta_0) + E_g[I_{D_2}(x) \nabla z_1^{(u)}(\theta_0|x)] = 0$ ,  
 (vii)  $V = (v^2 H(\theta_0) + E_g[I_{D_2}(x) \nabla^2 z_1^{(u)}(\theta_0|x)])^{-1}$  exists.

Then, (2)  $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_3(0, VWV')$  as  $n \rightarrow \infty$ , where  $W = (w_{uv})_{3 \times 3}$

is defined by

$$w_{uv} = q_1 \left( \frac{\partial \log(p_1(\theta))}{\partial \theta_u} \Big|_{\theta_0} \right) \left( \frac{\partial \log(p_1(\theta))}{\partial \theta_v} \Big|_{\theta_0} \right) + q_3 \left( \frac{\partial \log(p_3(\theta))}{\partial \theta_u} \Big|_{\theta_0} \right) \left( \frac{\partial \log(p_3(\theta))}{\partial \theta_v} \Big|_{\theta_0} \right) + E_g \left[ I_{D_2}(x) \left( \frac{\partial z_1^{(u)}(\theta|x)}{\partial \theta_u} \Big|_{\theta_0} \right) \left( \frac{\partial z_1^{(v)}(\theta|x)}{\partial \theta_v} \Big|_{\theta_0} \right) \right].$$

Although we have stated these results for Type I censoring, those for Type II are similar.

Figure 1

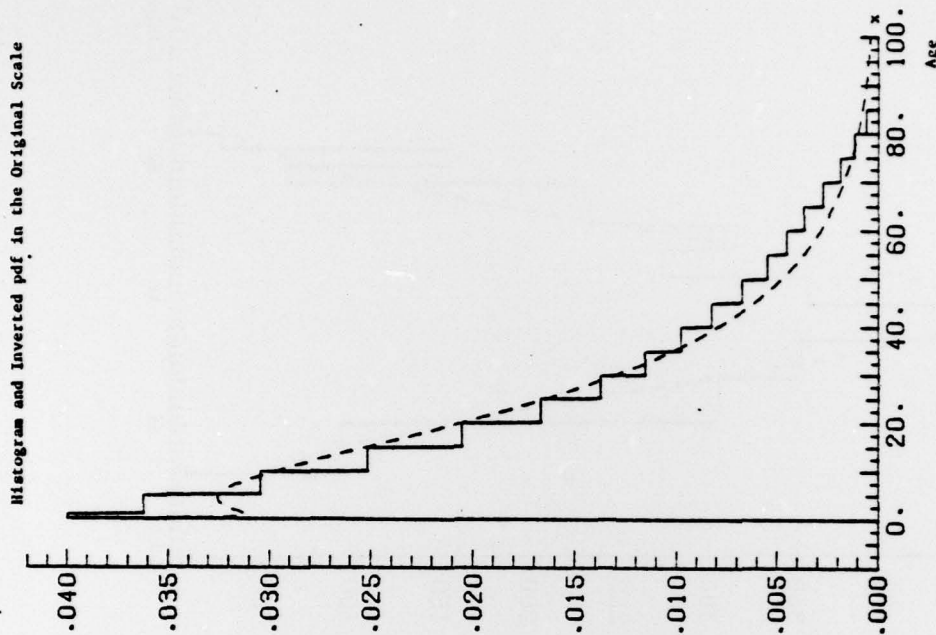


Figure 2

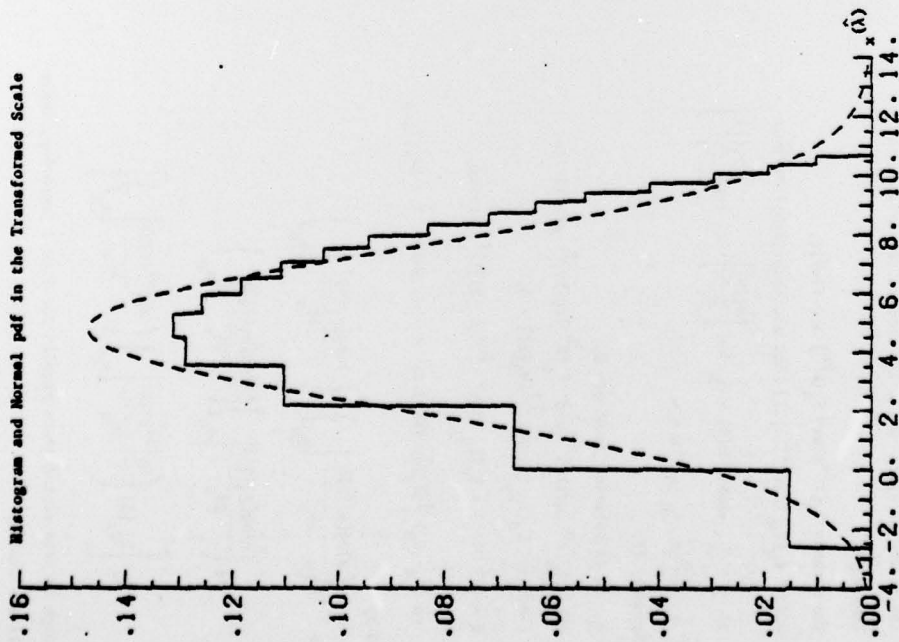
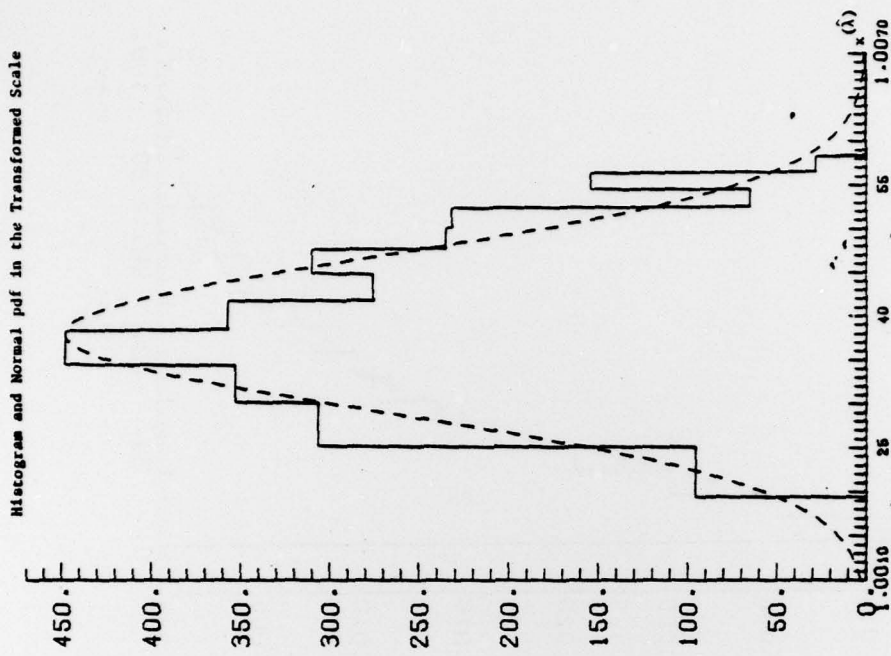
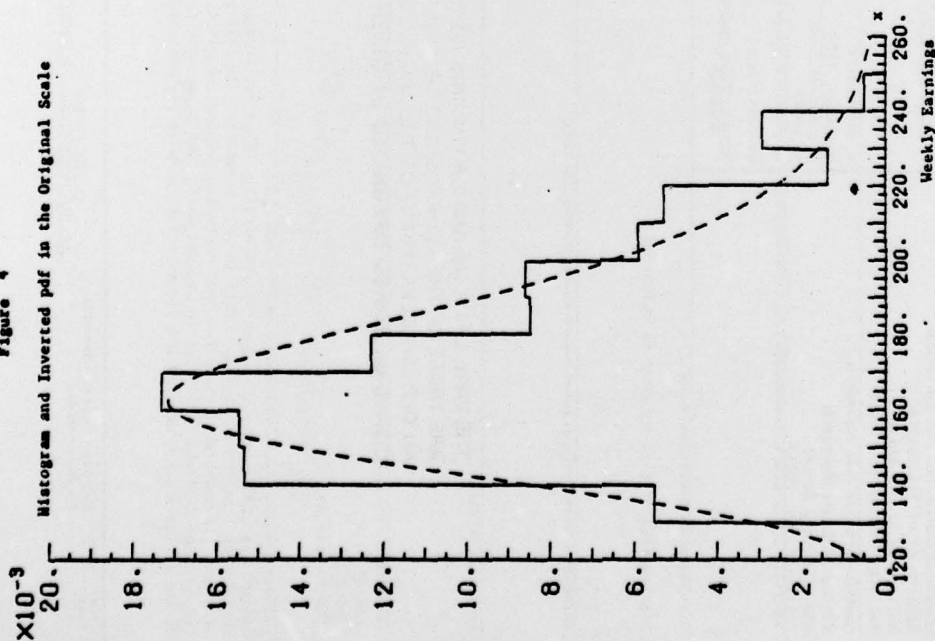


Figure 3



81b1lography

Figure 4



- Benn, R.T. and Sidebottom, S. (1976). "Algorithm AS-95: Maximum likelihood estimation of location and scale parameters from grouped data." Appl. Statist. 25, 88-93.
- Box, G.E.P. and Cox, D.R. (1964). "An analysis of transformations." J.R. Statist. Soc. B-26, 211-52.
- Draper, N.R. and Cox, D.R. (1969). "On distributions and their transformation to normality." J. Royal Statist. Soc. B-31, 472-76.
- Guerrero, G.V.M. (1979). "Extensions of the Box-Cox Transformation to Grouped-Data Situations." Ph.D. Thesis, University of Wisconsin-Madison.
- Hernández, A.F. and Johnson, R.A. (1979). "Transformation of a discrete distribution to near normality." TR No. 546, Dept. of Statistics, U. of Wisconsin-Madison.
- Kale, K.B. (1963). "Some remarks on a method of maximum-likelihood estimation proposed by Richards." J. Royal Statist. Soc. B-25, 209-12.
- Keyfitz, N. and Fliedger, M. (1971). Population, Facts and Methods of Demography. W.H. Freeman and Co.
- Lindley, D.V. (1950). "Grouping corrections and maximum likelihood equations." Proc. Camb. Philos. Soc. 46, 106-10.
- Richards, F.S.G. (1961). "A method of maximum likelihood estimation." J. Royal Statist. Soc. B-23, 469-75.

Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DDC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By	
Distribution/	
Availability Codes	
Dist	Avail and/or special
A	



Unclassified  
SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER Technical Report No. 542	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) TRANSFORMATION OF GROUPED OR CENSORED DATA TO NEAR NORMALITY		5. TYPE OF REPORT & PERIOD COVERED
6. AUTHOR(s) Victor M. Guerrero Richard A. Johnson		7. PERFORMING ORG. REPORT NUMBER
8. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Statistics University of Wisconsin Madison, Wisconsin 53706		9. CONTRACT OR GRANT NUMBER(s) ONR Grant No. N00014-78-C-0722 (Also funded by Army Res. Off)
10. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research 800 M. Quincy Street Arlington, VA 22217		11. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
12. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		13. REPORT DATE November, 1979
14. DISTRIBUTION STATEMENT (of this Report) Distribution of this document is unlimited		15. NUMBER OF PAGES 18 pages
16. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		17. SECURITY CLASS. (of this report) Unclassified
18. SUPPLEMENTARY NOTES		19. DECLASSIFICATION/DOWNGRADING SCHEDULE
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Grouped Data Sensory Transformations <i>An analysis of transformations</i>		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) Box and Cox (1964) proposed a power transformation which has proven utility for transforming ungrouped data to near normality. In this paper, we extend its applicability to two frequently encountered situations: (i) grouped data and (ii) Type I censored data. The asymptotic properties of the estimators are derived and illustrative examples are presented.		

DD FORM 1 JAN 73 1473 EDITION OF 1 NOV 65 IS OBSOLETE  
S/N 0102-LF-914-6601  
Unclassified  
SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

THE VIEW, OPINIONS, AND/OR FINDINGS CONTAINED IN THIS REPORT  
ARE THOSE OF THE AUTHOR(S) AND SHOULD NOT BE CONSTRUED AS  
AN OFFICIAL DEPARTMENT OF THE ARMY POSITION, POLICY, OR DE-  
CISION. UNLESS SO DESIGNATED BY OTHER DOCUMENTATION.